

Less is More: Learning to Refine Dialogue History for Personalized Dialogue Generation

Anonymous ACL submission

Abstract

Personalized dialogue systems explore the problem of generating responses that are consistent with the user’s personality, which have raised much attention in recent years. Existing personalized dialogue systems have tried to extract user profiles from dialogue history to guide personalized response generation. Since the dialogue history is usually long and noisy, most existing methods truncate the dialogue history to model the user personality. Such methods can generate some personalized responses, but a large part of dialogue history is wasted, leading to sub-optimal performance of personalized response generation. In this work, we propose to refine the user dialogue history from a large scale, based on which we can handle more dialogue history and obtain a more abundant and accurate persona information. Specifically, we design an MSP model which consists of three personal information refiners and a personalized response generator. With these multi-level refiners, we can sparsely extract the most valuable information (tokens) from the dialogue history and leverage other similar users’ data to enhance the personalization. Experimental results on two real-world datasets demonstrate the superiority of our model in generating more informative and personalized responses.¹

1 Introduction

Recent years have witnessed great progress on building personalized dialogue systems. In general, previous work explores building a personalized dialogue system via two pathways: (1) directly modelling user personality from predefined persona descriptions or user attribute (Qian et al., 2018; Zhang et al., 2018; Song et al., 2019); and (2) implicitly modelling the user personality from the user’s dialogue history (Li et al., 2016c; Ma et al., 2021). The latter is considered superior as the dialogue history is easy to obtain and comprises rich

personalized information. In this paper, we follow the second pathway that automatically learns implicit user profile from the user’s dialogue history to assist personalized response generation.

It is challenging to model user personality directly from the dialogue history. The reason is that a user’s dialogue history might contain massive historical dialogues, which are too heavy to load in the model and likely to be noisy. A straightforward solution is to truncate the dialogue history, as has been done by existing work (Ma et al., 2021; Qian et al., 2021a). However, as tremendous information has been wasted, the model’s performance is also influenced. On the other hand, we observe that the dialogue history from other users’ may also be helpful in generating a personalized response for the current user. For example, users with the same interest on “soccer” may talk about similar things on such a topic. This has been overlooked by existing methods. Intuitively, the problem of “data explosion” is even more severe in the latter case (when other similar users’ dialogue history is also considered). To alleviate these problems, we propose using a hierarchical refiner structure to sparsely extract the most valuable query-aware persona information from both the current and other similar users’ dialogue history. By this means, more dialogue history can be utilized for learning user personality and improving response generation.

Our model is called **MSP**, which stands for **M**odeling and **S**electing user **P**ersonality from the dialogue history for generating personalized responses. Instead of attending to all dialogue history, MSP refines the most valuable historical information that can well portray the user personality and guide the response generation. Specifically, MSP consists of three personal information refiners working in different levels and a personalized response generator. At **first**, a user refiner is designed to select a group of users who have similar interests with the current user. By refining dialogue

¹Our codes are released in anonymous.4open.science.

history at the user level, we can obtain similar data to share information with similar users and avoid mutual interference with other users. **Then**, a topic refiner filters out the current and similar users’ dialogue history that has different topics with the current query in sentence level. **Next**, we design a token refiner to extract the most valuable query-aware user profiles from the remaining dialogue history in the token level. **Finally**, a personalized response generator combines user profiles and the current query to generate responses. Given that there is no explicit supervisory signal guiding the refiner extract an exemplary user profile, we design a supplementary sentence matching task and a joint training method. The generator will construct a pseudo-label to guide the refiner’s extraction.

Our contributions are three-fold: (1) We design an MSP model to tackle the data noise problem. It can efficiently refine user profiles through dialogue history and generate personalized responses. (2) We design a refiner structure to extract the query-aware profile in three levels. The similar users’ information is taken into account, which can help improve the personality for the response. (3) We design a joint training method of the refiner and generator. The refiner provides the generator with user profiles to assist in generating responses, while the generator constructs a pseudo-label for the refiner to assist in selecting user profiles.

2 Related Work

Personalized dialogue generation. Open-domain dialogue generation has been extensively studied (Koehn et al., 2003; Vinyals and Le, 2015; Serban et al., 2016; Zhang et al., 2019a,b; Xiao et al., 2020). Recently, personalized dialogue systems have attracted more and more attention. Typical methods include: (1) explicitly using predefined persona descriptions or attributes as users’ profile to generate personalized responses (Qian et al., 2018; Zhang et al., 2018; Olabiyi et al., 2019; Song et al., 2019); (2) using user ID embeddings to enhance personalized dialogue generation (Li et al., 2016c; Chan et al., 2019); and (3) extracting implicit user profile from users’ dialogue history to generate personalized responses (Al-Rfou et al., 2016; Bak and Oh, 2019). Since manually collecting user profiles is impractical for large-scale datasets and the user ID embeddings perform badly, in this study, we focus on the last group of methods for personalized response generation.

DHAP (Ma et al., 2021) is the state-of-the-art method in personalized dialogue generation. It uses a transformer-based structure to model the user’s dialogue history and extract personal information for response generation. Unfortunately, this model can only handle a limited number of user dialogue histories, wasting a lot of valuable information. Our method has two main differences with DHAP: (1) We propose a refiner structure in our model so that more dialogue history can be handled and the most valuable information can be extracted for improving response generation. (2) With our proposed refiner, we can further incorporate more dialogue history from other users (having similar interests) to facilitate personalized dialogue generation.

Retrieval-guided natural language generation. Retrieval-based methods can collect relevant information for language generation (Yang et al., 2019). It has been widely applied in many tasks such as text style transfer (Li et al., 2018) and dialogue generation (Wu et al., 2019; Cai et al., 2019). The idea of using a retrieval system to get useful information inspires our study. We use a refiner to automatically extract personal information from dialogue history and guide the personalized generation.

3 Methodology

3.1 Problem Statement and Overview

Considering a set of users $\mathcal{U} = \{u_1, \dots, u_l\}$, for any user u_i , we have their dialogue history with others $U^i = \{(q_1^i, r_1^i), \dots, (q_t^i, r_t^i)\}$, where q_j^i is a *query* issued by others, while r_j^i is the corresponding *response* given by u_i .² Our target is to generate a personalized response r^i for the user u_i to reply a new query q . As we introduced earlier, the personalized information can be obtained from the dialogue history U^i of the user u_i and other dialogue history U^j from similar users $u_j (j \neq i)$.

The overview of our MSP model is shown in Figure 1. MSP consists of three personal information refiners working in different levels and a personalized response generator. Specifically, the first refiner is working at the user-level. By comparing the dialogue history of the current user u_i with others, MSP can select a group of users having similar interests with u_i . After obtaining a group of similar users, we further refine their dialogue

²Here we use the term “query” to denote the utterance given by others. Generally, the query can be either one utterance in single-turn dialogues, or several history utterances in multi-turn dialogues.

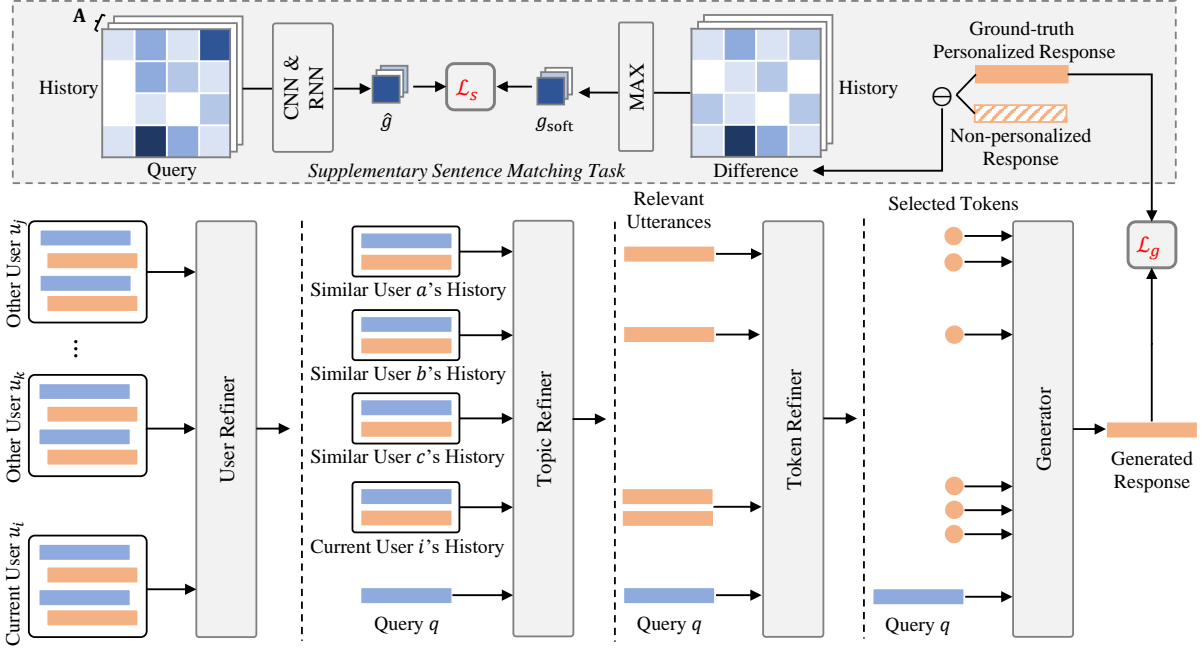


Figure 1: The overview structure of the proposed model which consists of four modules: (1) user refiner, (2) topic refiner, (3) token refiner, and (4) generator.

history according to the relevance with the current query’s topic. Moreover, we add the last refiner to extract several tokens so that the most fine-grained personal information can be extracted from the relevant utterances. Finally, the query and the extracted tokens are fed into the generator and construct a personalized response.

3.2 User Refiner

The dialogue history of users with similar interests may share much personal information. Therefore, our first target is to select a group of users with similar interests to the current user. We design a user refiner to achieve this. Since the users’ interest is usually contained in their dialogues with others, we consider both the queries and responses in the dialogue history to select similar users. Specifically, for the user u_i ’s dialogue history U^i , we apply a pre-trained BERT (Devlin et al., 2019) and represent them by the embedding of “[CLS]” token:

$$\mathbf{U}_q^i = \sum_{j=1}^t \text{BERT}(q_j^i), \quad \mathbf{U}_r^i = \sum_{j=1}^t \text{BERT}(r_j^i).$$

Then, we can select k_u users that have similar interest with the current user u_i :

$$u^{\text{sim}} = \text{TopK}(\mathbf{U}^i \cdot \mathbf{U}^j, k_u), \quad (1)$$

$$\mathbf{U}^i = [\mathbf{U}_q^i; \mathbf{U}_r^i], \quad (2)$$

where $\text{TopK}(\cdot, \cdot)$ is the top- k selection operation.

After the user refiner, we can obtain the dialogue history of the similar users $\{u_j\}_{j=1}^{k_u}$. It is worth noting that, since the number of users is large in the datasets, we choose to use the dot-product to compute the similarity of the users so that the whole process can be implemented by dense retrieval libraries, such as Faiss (Johnson et al., 2017), which is very efficient.

3.3 Topic Refiner

The users’ dialogue history often contains many dialogues with others. These dialogues have various topics, which may be irrelevant to the current query. Therefore, we propose a topic refiner to select relevant dialogue history for personalized response generation. Specifically, we use a topic classifier to compute the topic distribution of the current query q and the queries q_j^i in the history dialogues:

$$t = \text{MLP}(\text{mean}(\text{BERT}(q))), \quad (3)$$

$$t_j^i = \text{MLP}(\text{mean}(\text{BERT}(q_j^i))), \quad (4)$$

where $t, t_j^i \in \mathbb{R}^{d^t \times 1}$, and d^t is the number of topic. Then, we filter out the dialogue history $\langle q_j^i, r_j^i \rangle$ that has different topics with the current query, i.e., $\max(t_j^i) \neq \max(t)$.

In the topic refining process, we compare the queries in the history dialogues with the current query to filter out topic-irrelevant dialogues. This

process can further reduce the noise and make our model more lightweight. Both the dialogue history of the current user and that of the similar users (obtained in the former step) are refined. In the next step, we will use the **responses** in these selected history dialogues and extract the most valuable tokens for personalized response generation.

3.4 Token Refiner

After the previous two refiners, we obtain a collection of historical responses. Though we can directly add them into the generation process, our preliminary experiments indicate that they perform poorly. A major reason is the noisy quality of the responses. Indeed, existing studies (Shang et al., 2015; Borgeaud et al., 2021) have demonstrated the effectiveness of using informative tokens to improve the response generation. Inspired by these studies, we further devise a token refiner to extract the most fine-grained information (tokens) from the historical responses. Specifically, we compute an attention map \mathbf{A} between the query q and the historical responses r^{sim} and r^{cur} (they are from the similar users and the current user respectively) as:

$$\mathbf{A} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right), \quad (5)$$

$$\mathbf{Q} = \text{TRM}_{\text{enc}}(q) \cdot \mathbf{W}_Q, \quad (6)$$

$$\mathbf{K} = \text{TRM}_{\text{enc}}(r) \cdot \mathbf{W}_K, \quad (7)$$

where $\text{TRM}_{\text{enc}}(\cdot)$ is a transformer encoder. r refers to r^{sim} or r^{cur} , and correspondingly, \mathbf{A} refers to the similar user matching map \mathbf{A}^{sim} and current user matching map \mathbf{A}^{cur} . $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times d}$ are parameters, and d is the dimension of the hidden state. After obtaining the attention matching map \mathbf{A} , we select tokens to form the similar users' profile and current user's profile according to each token's attention weight:

$$\mathbf{c}^{\text{sim}} = \text{TopK}(\text{Max}(\mathbf{A}^{\text{sim}}), k_p), \quad (8)$$

$$\mathbf{c}^{\text{per}} = \text{TopK}(\text{Max}(\mathbf{A}^{\text{cur}}), k_p), \quad (9)$$

where k_p is a hyper-parameter to control the number of profile tokens.

3.5 Generator

We use a transformer decoder as to generate a personalized response by using the similar users' profile \mathbf{c}^{sim} , current user's profile \mathbf{c}^{cur} , and query information \mathbf{q} as input. The decoding process can be

defined as:

$$\hat{y} = \text{TRM}_{\text{dec}}(\mathbf{x}), \quad (10)$$

$$\mathbf{x} = [\mathbf{c}^{\text{sim}}; \mathbf{c}^{\text{per}}; \mathbf{q}], \quad (11)$$

where $[\cdot]$ is the concatenation operation and \hat{y} is the word generation probability.

3.6 Training and Optimization

The generator is optimized by maximizing the generation probability of the ground-truth y :

$$\mathcal{L}^g = -y \log \hat{y}. \quad (12)$$

In practice, we find that the token refiner is hard to train. We speculate the reason is a missing of direct supervision signals. In this case, it is difficult to tell whether the training errors stem from the generation process or the refining process. To tackle this problem, we propose a supplementary sentence matching task to assist the token selection.

Supplementary sentence matching task. The core idea of this task is to train the token refiner directly by introducing supervision signals so that it can automatically mine valuable tokens. Specifically, we design a sentence matching task to match the query with the dialogue history. The task's target is to find the history sentences that help generate personalized responses. We consider using the query-history cross attention weight \mathbf{A} to generate a matching representation and then use this representation to finish the task. In this way, once the matching task has been well-finished, we can use attention map \mathbf{A} to identify the most informative and valuable parts of a history sentence that are helpful to generate a personalized response.

To achieve our idea, we design a matching process. Firstly, we calculate the matching representations \mathbf{H} by the cross-attention map \mathbf{A} and then apply a CNN with a max-pooling operation to aggregate token information:

$$\mathbf{S} = \text{Maxpool}(\text{CNN}(\mathbf{H})) \quad (13)$$

$$\mathbf{H} = \mathbf{A} \cdot \mathbf{V}, \quad (14)$$

$$\mathbf{V} = \text{TRM}_{\text{enc}}(r) \cdot \mathbf{W}_V. \quad (15)$$

Next we flatten \mathbf{S} and applies a LSTM to aggregate the sentence information:

$$\mathbf{h} = \text{LSTM}(\text{Flatten}(\mathbf{S})). \quad (16)$$

Finally we use sentence matching vector \mathbf{h} to compute the matching score:

$$\hat{g} = \text{Sigmoid}(\text{MLP}(\mathbf{h})). \quad (17)$$

Algorithm 1 Joint training process

Input: M dialogue triplets: $D = \{ \langle q_i, y_i, r_i \rangle \}_{i=1}^M$ **Output:** A personalized dialogue model

```
1: Init the refiner and generator module
2: while not converge do
3:   Sample  $n_s$  dialogue triplets  $D^q = \{ \langle q_i, y_i, r_i \rangle \}_{i=1}^{n_s}$ 
4:   Get  $\hat{Y}^q = \{ \hat{y}_i \}_{i=1}^{n_s}$  on  $D^q$  from  $p_g(\hat{y}|q)$ 
5:   Get pseudo-label  $g$  on  $D^q$ 
6:   Train refiner by optimizing  $\mathcal{L}^s$  on  $D^q$ 
7:   if Current Step  $> N_f$  then
8:     Sample  $n_d$  dialogue triplets  $D^p = \{ \langle q_i, y_i, r_i \rangle \}_{i=1}^{n_d}$ 
9:     Extract  $C^p = \{ c_i \}_{i=1}^{n_d}$  on  $D^p$  from  $p_s(c|q, r)$ 
10:    Train generator by optimizing  $\mathcal{L}^g$  on  $D^p \cup C^p$ 
11:    end if
12: end while
```

For guiding the sentence matching task, we design a pseudo-label g to measure the matching goodness of each history sentence. We expect the history with more persona profile information can achieve a higher score. Thus we use the difference between personalized ground-truth y and non-personalized generated response probability \hat{y} to measure the persona profile and create the pseudo-label:

$$g = \begin{cases} 1, & g_{soft} \geq \alpha \\ 0, & g_{soft} < \alpha, \end{cases} \quad (18)$$

$$g_{soft} = \text{Sum}(\text{Max}((y - \hat{y}) \cdot r)) / d_y, \quad (19)$$

where d_y is the length of ground-truth y and α is a threshold. Finally, we minimize the binary cross entropy loss between g and \hat{g} :

$$\mathcal{L}^s = g \log \hat{g} + (1 - g) \log(1 - \hat{g}). \quad (20)$$

Joint training. To facilitate the learning with the above gradient approximation approach, we design a joint training process to train the refiner and generator in turn. Specifically, in each training iteration, we first sample a batch of query q , response y , similar and current users' dialogue history r^{sim} and r^{per} from dataset D . After through a non-personalized generator, we create the pseudo-label g (Eq.18). This pseudo-label is used to train the token refiner by optimizing the loss \mathcal{L}^s (Eq. 20). Further, we sample another batch D^p from D . After extracting sim profile c^{sim} and persona profile c^{per} (Eq.8, Eq.9), we generate the personalized response \hat{y} and update the generator by optimizing the loss \mathcal{L}^g (Eq. 12). To avoid bad profile misleading the generation at the beginning training process, we pre-train the refiner for N_f steps before extracting the profile for the generator. The detailed training process is summarized in Algorithm 1.

4 Experiments

4.1 Datasets

To evaluate our model's performance, we conduct experiments on a Chinese Weibo³ dataset (Qian et al., 2021b) and an English Reddit⁴ dataset. Both are collected from open-domain social media platforms. On these platforms, users can post various topics, and other users can respond to them. We compare user-id and timestamps to associate the query with its corresponding response and the current user's dialogue history. As a result, each training sample contains a query, a response, and a sequence of dialogue history. Finally, the dataset is divided into training, validation, and test sets in chronological order. The statistics of the datasets are provided in Appendix A.

4.2 Baselines

We compare our proposed model with eight highly correlated and strong baselines.⁵ They can be categorized into four groups:

Non-personalized methods. (1) Seq2Seq-Attention (Sutskever et al., 2014) is a vanilla sequence-to-sequence model with attention mechanism (Luong et al., 2015). (2) MMI (Li et al., 2016a) is based on seq2seq and use maximum mutual information as an extra loss to improve diversity. (3) DialoGPT (Zhang et al., 2019b) is a variant of GPT-2 (Radford et al., 2019) designed for dialogue generation.

Predefined profile-based methods. Since there are no persona descriptions in the datasets, we test these methods by using the user's dialogue history as a simulation of predefined persona profiles. (4) GPMN (Zhang et al., 2018) enhances the seq2seq model with a memory module, which encodes and stores the persona profile as memory representations. (5) PerCVAE (Zhao et al., 2017) encodes predefined personalized sentences as a conditional representation and uses CVAE to generate a personalized response.

User ID-based methods. (6) Speaker (Li et al., 2016c) is based on seq2seq while using user-id embedding as user representation to facilitate the response generation. (7) Persona WAE (Chan et al., 2019) uses WAE (Wasserstein autoencoder) for response generation. It maps user-id embeddings to a personalized Gaussian mixture distribution and

³<https://www.weibo.com/>

⁴<https://www.reddit.com/>

⁵The implementation details are given in Appendix B.

399 then samples the personalized vector to guide the
400 response generation.

401 **User dialogue history-based methods.** (8)
402 DHAP (Ma et al., 2021) uses history memory to
403 store and construct the dynamic query-aware user
404 profile from dialogue history and then uses a per-
405 sonalized decoder to generate a response. Since
406 this model also learns the user profile directly from
407 the dialogue history, it is the most relevant baseline
408 of our method.

409 4.3 Evaluation

410 **Metric-based.** We first evaluate all methods by
411 several metrics with respect to different aspects.
412 (1) BLEU-1/2 (Papineni et al., 2002) and ROUGE-
413 L (Lin and Och, 2004) are typical word overlap-
414 based metrics for measuring the similarity between
415 the generated response and the ground-truth.⁶ (2)
416 Distinct-1/2 (Li et al., 2016b) consider the num-
417 ber of uni- or bi-grams in the generated response,
418 which is commonly used for evaluating the diver-
419 sity. (3) The embedding-based metrics (*i.e.*, aver-
420 age, extrema, and greedy) (Chan et al., 2019) are
421 introduced to measure the semantic similarity be-
422 tween the generated response and the ground-truth
423 one. (4) As a personalized dialogue model, follow-
424 ing previous studies (Ma et al., 2021), two tailored
425 metrics are adopted to measure how much informa-
426 tion is included in the dialogue history that can be
427 reflected in the response. Persona-F1 (P-F1) (Lian
428 et al., 2019) calculates the F1 value to measure
429 the uni-grams co-occurring in both the generated
430 response and the dialogue history. Persona Cov-
431 erage (P-Cover) (Song et al., 2019) calculates the
432 IDF-weighted word overlap between the generated
433 response and the golden one so that the importance
434 of different words can be taken into account.

435 **Human Annotation.** Due to the variability of
436 human language, a response that differs from the
437 ground-truth may also be appropriate. Following
438 previous studies (Chan et al., 2019), we conduct a
439 human evaluation of all methods. Concretely, we
440 sample 100 (query, response, user dialogue history)
441 triplets and hire three well-educated annotators to
442 score the responses generated by different mod-
443 els. Three aspects, *i.e.*, readability, informativ-
444 eness, and personalization, are considered. The first
445 two factors are scored on a scale of [1, 3] for their
446 quality, while the third is assessed on a scale of [0,

1], indicating whether the response can accurately
447 reflect the user’s personality.⁷ 448

449 4.4 Experimental Results

450 **Metric-based evaluation.** Table 1 shows all mod-
451 els’ performance under different metrics. On both
452 datasets, it is clear to see that our MSP model out-
453 performs baselines on all metrics. The improve-
454 ment is statistically significant (t-test with p -value
455 < 0.05). These findings indicate that our model
456 is capable of generating more fluent, diverse, and
457 personalized responses than all baselines. In par-
458 ticular, we can observe: (1) MSP achieves better
459 performance on overlap-based metrics. This sug-
460 gests that our model can provide responses that are
461 more similar to the ground-truth with the help of
462 the selected tokens. (2) For diversity metrics, the
463 higher distinct values show that our generated re-
464 sponses are more diverse. Additionally, predefined
465 profile-based methods and user dialogue history-
466 based methods outperform others. This shows that
467 incorporating external information can aid in the
468 generation of more informative responses. (3) In
469 addition to generating more overlapped words with
470 the ground-truth response, the improvements of
471 embedding metrics reflect that our model generates
472 more semantically relevant responses. (4) Finally,
473 the increase in personalized metrics implies that
474 our approach can incorporate more user-specific
475 information into the generation. Furthermore, the
476 significant improvement over DHAP demonstrates
477 that our model can extract more meaningful person-
478 alized information from the user dialogue history.

479 **Human annotation.** The result of human anno-
480 tation on the Weibo dataset is shown in Table 2.
481 The Fleiss Kappa is around 0.62, indicating a sub-
482 stantial agreement achieved by three annotators. In
483 general, the results of human annotation are con-
484 sistent with those of the metric-based evaluation.
485 Both of them demonstrate our model’s superiority
486 at generating more fluent, informative, and person-
487 alized responses. Compared to non-personalized
488 methods, user id-based methods can enhance per-
489 sonalization at the expense of readability. User
490 dialogue history-based methods (*i.e.*, DHAP and
491 MSP) can largely improve the personalization of
492 the response while retaining a high level of read-
493 ability and informativeness. We attribute this to
494 the abundant personal information contained in the
495 user dialogue history.

⁶The results of BLEU-3/4 and ROUGE-1/2 are provided
in Appendix C.

⁷The detailed scoring criteria are described in Appendix D.

Table 1: The result of metric-based evaluation on the Weibo dataset and Reddit Dataset. “†” indicates that our model achieves significant improvement in t-test with p -value < 0.05 .

	Overlap-based Metric			Diversity		Embedding Metric			Persona Metric		
	BLEU-1	BLEU-2	ROUGE-L	Dist-1	Dist-2	Average	Extrema	Greedy	P-F1	P-Cover	
Weibo	Seq2Seq	3.330†	0.294†	8.985†	0.930†	2.180†	0.321†	0.266†	0.254†	0.154†	0.041†
	MMI	3.631†	0.095†	5.264†	10.710†	43.458†	0.477†	0.696†	0.305†	0.325†	0.054†
	DialoGPT	6.068†	0.741†	8.459†	15.322†	55.536†	0.557†	0.793†	0.324†	0.522†	0.061†
	GPMN	4.899†	0.696†	7.785†	11.724†	32.730†	0.353†	0.391†	0.301†	0.542†	0.084†
	PerCVAE	5.114†	0.299†	7.380†	14.098†	49.733†	0.469†	0.657†	0.299†	0.903†	0.086†
	Speaker	4.994†	0.113†	7.868†	6.035†	19.007†	0.492†	0.712†	0.311†	0.225†	0.082†
	PersonaWAE	3.510†	0.155†	10.546†	2.551†	19.743†	0.563†	0.757†	0.307†	1.740†	0.103†
	DHAP	9.324†	0.894†	14.122†	15.175†	58.806†	0.523†	0.747†	0.313†	1.791†	0.144†
	MSP (Ours)	11.875	5.108	15.563	24.203	73.196	0.605	0.883	0.331	2.170	0.297
Reddit	Seq2Seq	1.820†	0.023†	4.069†	5.203†	19.485†	0.545†	0.554†	0.470†	0.051†	0.029†
	MMI	2.065†	0.011†	3.784†	5.914†	31.093†	0.543†	0.607†	0.454†	0.085†	0.038†
	DialoGPT	4.735†	0.397†	8.943†	6.353†	29.106†	0.604†	0.733†	0.448†	0.137†	0.040†
	GPMN	2.686†	0.376†	4.776†	12.325†	35.762†	0.406†	0.331†	0.358†	0.189†	0.037†
	PerCVAE	5.933†	0.576†	8.112†	9.631†	40.213†	0.637†	0.649†	0.499†	0.212†	0.040†
	Speaker	2.642†	0.054†	4.469†	8.951†	34.187†	0.538†	0.606†	0.457†	0.115†	0.031†
	PersonaWAE	2.637†	0.113†	8.199†	1.758†	25.915†	0.629†	0.685†	0.442†	0.206†	0.032†
	DHAP	6.858†	0.737†	11.720†	18.707†	66.932	0.709	0.721†	0.539	0.227†	0.111†
	MSP (Ours)	7.174	0.883	12.171	21.247	68.897	0.716	0.764	0.545	0.276	0.137

Table 2: The result of human evaluation on Weibo dataset. “†” indicates that our model achieves significant improvement in t-test with p -value < 0.05 .

Model	Readability	Informativeness	Personality
Seq2Seq	1.76†	1.37†	0.11†
MMI	1.96†	1.88†	0.19†
DialoGPT	2.33	2.10†	0.32†
GPMN	2.01†	2.16†	0.35†
PerCVAE	2.10†	2.01†	0.39†
Speaker	1.89†	1.44†	0.24†
PersonaWAE	1.81†	2.01†	0.32†
DHAP	2.29†	2.19†	0.55†
MSP (Ours)	2.37	2.39	0.67
Ground-Truth	2.71	2.66	0.76

Table 3: The results of ablation experiments on Weibo dataset.

Models	BLEU-1	BLEU-2	P-Cover
MSP (Full)	11.875	5.108	0.297
w/o User Refiner	6.093	0.757	0.151
w/o Topic Refiner	6.163	0.839	0.178
w/o Token Refiner	4.213	0.609	0.116
w/o Current U’s Profile	9.365	3.146	0.238
w/o Similar Us’ Profile	6.413	0.871	0.245
w/o Joint Training	6.070	0.749	0.130

5 Further Analysis

We further conduct a series of analyses to elaborate our model. All analysis here is based on the Weibo dataset, while similar results can be observed on the Reddit dataset.

Ablation Study. To investigate the impact of different modules in MSP, we conduct an ablation study by removing or using different strategies in each module.

We first study the influence of the refiners at three levels. (1) We remove the user refiner and train our model using randomly sampled users. We can see the performance of all metrics drops. This illustrates that our MSP model can select users that shares the same interests as the current user and thereby improving response quality. (2) We remove the topic refiner and supply the token refiner with full dialogue history. The performance degradation demonstrates that various topics in dialogue history introduce lots of noise, misleading the token refiner on extracting valuable tokens, thus impairing the personalized response generation. (3) We eliminate the token refiner and feed all dialogue history sentences directly into the generator.⁸ The decline in performance implies the effectiveness and necessity of token selection. It is worth noting that, as compared to using the complete history, our selection strategy can reduce training time by 41.6%, considerably increasing efficiency. All of the aforementioned experimental results suggest that MSP’s advantage stems from high-quality personalized information extraction rather than simply introducing additional information.

⁸Due to the length limitation of GPT-2, history with more than 512 tokens will be truncated.

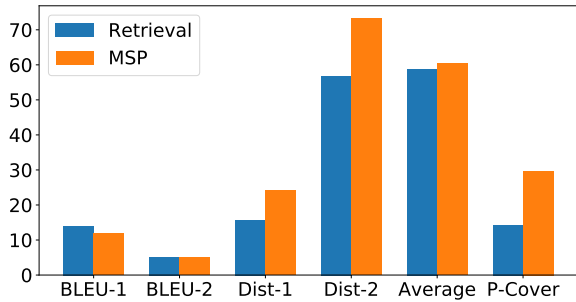


Figure 2: Comparison with the retrieval-based model on the Weibo dataset.⁹

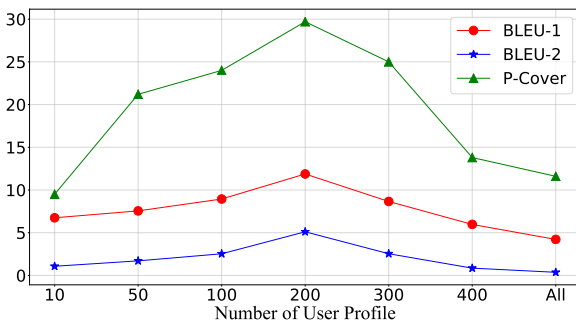


Figure 3: Experiments with the different number of user profiles on the Weibo dataset.⁹

We then explore the impact of personalized information from two sources, *i.e.*, the current user’s profile and the similar users’ profile. Removing either of them results in decreased performance. This exemplifies their usefulness. Specifically, compared with similar users’ profiles, eliminating the current user’s profile will hurt the personalization effect heavily. This result shows that, for personalization, the current user’s profile is more essential than that of similar users, which is quite intuitive. However, the similar users’ profile has a significant effect on BLEU-1/2, implying that such profile can provide abundant information for response generation. Consequently, integrating both types of profiles significantly improves response generation.

Finally, we conduct an experiment to validate our proposed joint training for the token refiner. The declining performance indicates that the token refiner is unable to extract useful information in the absence of additional supervision signals. Indeed, when the sentence matching task is removed, the token refiner extracts tokens that are relevant to the current query, which is less useful for generating a personalized response.

Influence of the selection mechanism. To val-

⁹To keep the dimensionality consistent, P-Cover is multiplied by a factor of 100.

idate the effectiveness of our proposed selection mechanism, we replace the refiner with a traditional retrieval method (*i.e.*, BM25 (Robertson and Walker, 1994)). Specifically, we use the query to retrieve 15 relevant responses and feed them into our model for training. The experimental results are shown in Figure 2. We can observe that the retrieval strategy achieves comparable performance with our model on word-overlap and embedding-based metrics. This suggests that the relevant dialogue history for the query can provide valuable information for response generation. However, the retrieval strategy performs poorly on diversity and personalization metrics. This demonstrates that, without careful selection, the retrieved information is too generic and thus less helpful for personalized response generation.

Influence of the personalized tokens amount. In MSP, three refiners are designed to extract personalized tokens for response generation. Intuitively, the amount of the tokens will have an effect on the refiner’s performance. We report this influence in Figure 3. As we can see in the left part, the quality of response generation improves with more tokens used. This is because fewer tokens are incapable of covering sufficient personalized information for response generation. Our MSP model performs optimally when about 200 personalized tokens are selected. When more tokens are introduced, the performance degrades. The potential reason is that more tokens would bring noise to the generation. This is consistent with our speculation that the dialogue history is noisy and the information selection is both effective and necessary.¹⁰

6 Conclusion

In this work, we propose an MSP model for personalized response generation. Unlike previous related work, we utilize a refiner structure to extract query-aware persona information from large-scale dialogue history. The multi-level refiners can sparsely extract valuable information from dialogue history and leverage similar users’ information to enhance the current user’s personalization. Experimental results confirm the effectiveness of our model on generating informative and personalized responses.

¹⁰Due to the space limitation, we present a case study in Appendix E.

References

- 599
- 600 Rami Al-Rfou, Marc Pickett, Javier Snaider, Yun-Hsuan
601 Sung, Brian Strope, and Ray Kurzweil. 2016. [Con-](#)
602 [versational contextual cues: The case of personal-](#)
603 [ization and history for response ranking.](#) *CoRR*,
604 [abs/1606.00372](#).
- 605 JinYeong Bak and Alice Oh. 2019. [Variational hierarchi-](#)
606 [cal user-based conversation model.](#) In *Proceedings*
607 *of the 2019 Conference on Empirical Methods in Nat-*
608 *ural Language Processing and the 9th International*
609 *Joint Conference on Natural Language Processing,*
610 *EMNLP-IJCNLP 2019, Hong Kong, China, Novem-*
611 *ber 3-7, 2019*, pages 1941–1950. Association for
612 Computational Linguistics.
- 613 Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann,
614 Trevor Cai, Eliza Rutherford, Katie Millican, George
615 van den Driessche, Jean-Baptiste Lespiau, Bogdan
616 Damoc, Aidan Clark, Diego de Las Casas, Aurelia
617 Guy, Jacob Menick, Roman Ring, Tom Hennigan,
618 Saffron Huang, Loren Maggiore, Chris Jones, Albin
619 Cassirer, Andy Brock, Michela Paganini, Geoffrey
620 Irving, Oriol Vinyals, Simon Osindero, Karen Si-
621 monyan, Jack W. Rae, Erich Elsen, and Laurent Sifre.
622 2021. [Improving language models by retrieving from](#)
623 [trillions of tokens.](#) *CoRR*, [abs/2112.04426](#).
- 624 Deng Cai, Yan Wang, Wei Bi, Zhaopeng Tu, Xi-
625 aojiang Liu, and Shuming Shi. 2019. [Retrieval-](#)
626 [guided dialogue response generation via a matching-](#)
627 [to-generation framework.](#) In *Proceedings of the*
628 *2019 Conference on Empirical Methods in Natu-*
629 *ral Language Processing and the 9th International*
630 *Joint Conference on Natural Language Processing,*
631 *EMNLP-IJCNLP 2019, Hong Kong, China, Novem-*
632 *ber 3-7, 2019*, pages 1866–1875. Association for
633 Computational Linguistics.
- 634 Zhangming Chan, Juntao Li, Xiaopeng Yang, Xiuying
635 Chen, Wenpeng Hu, Dongyan Zhao, and Rui Yan.
636 2019. [Modeling personalization in continuous space](#)
637 [for response generation via augmented wasserstein](#)
638 [autoencoders.](#) In *Proceedings of the 2019 Confer-*
639 *ence on Empirical Methods in Natural Language Pro-*
640 *cessing and the 9th International Joint Conference*
641 *on Natural Language Processing, EMNLP-IJCNLP*
642 *2019, Hong Kong, China, November 3-7, 2019*, pages
643 1931–1940. Association for Computational Linguis-
644 tics.
- 645 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
646 Kristina Toutanova. 2019. [BERT: pre-training of](#)
647 [deep bidirectional transformers for language under-](#)
648 [standing.](#) In *Proceedings of the 2019 Conference of*
649 *the North American Chapter of the Association for*
650 *Computational Linguistics: Human Language Tech-*
651 *nologies, NAACL-HLT 2019, Minneapolis, MN, USA,*
652 *June 2-7, 2019, Volume 1 (Long and Short Papers),*
653 *pages 4171–4186.* Association for Computational
654 Linguistics.
- 655 Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and
656 Yejin Choi. 2020. [The curious case of neural text](#)
[degeneration.](#) In *8th International Conference on*
Learning Representations, ICLR 2020, Addis Ababa,
Ethiopia, April 26-30, 2020. OpenReview.net.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. Billion-scale similarity search with gpus. *arXiv preprint arXiv:1702.08734*.
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A](#)
[method for stochastic optimization.](#) In *3rd Inter-*
national Conference on Learning Representations,
ICLR 2015, San Diego, CA, USA, May 7-9, 2015,
Conference Track Proceedings.
- Philipp Koehn, Franz Josef Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology-Volume 1*, pages 48–54. Association for Computational Linguistics.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016a. [A diversity-promoting ob-](#)
[jective function for neural conversation models.](#) In *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, pages 110–119. The Association for Computational Linguistics.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016b. A diversity-promoting objective function for neural conversation models. In *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, pages 110–119.
- Jiwei Li, Michel Galley, Chris Brockett, Georgios P. Spithourakis, Jianfeng Gao, and William B. Dolan. 2016c. A persona-based neural conversation model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers.*
- Juncen Li, Robin Jia, He He, and Percy Liang. 2018. [Delete, retrieve, generate: a simple approach to senti-](#)
[ment and style transfer.](#) In *Proceedings of the 2018*
Conference of the North American Chapter of the
Association for Computational Linguistics: Human
Language Technologies, NAACL-HLT 2018, New Or-
leans, Louisiana, USA, June 1-6, 2018, Volume 1
(Long Papers), pages 1865–1874. Association for
Computational Linguistics.
- Rongzhong Lian, Min Xie, Fan Wang, Jinhua Peng, and Hua Wu. 2019. [Learning to select knowledge for re-](#)
[sponse generation in dialog systems.](#) In *Proceedings*
of the Twenty-Eighth International Joint Conference
on Artificial Intelligence, IJCAI 2019, Macao, China,
August 10-16, 2019, pages 5081–5087. ijcai.org.

825 Hainan Zhang, Yanyan Lan, Liang Pang, Jiafeng Guo,
826 and Xueqi Cheng. 2019a. [Recosa: Detecting the](#)
827 [relevant contexts with self-attention for multi-turn](#)
828 [dialogue generation](#). In *Proceedings of the 57th Con-*
829 *ference of the Association for Computational Lin-*
830 *guistics, ACL 2019, Florence, Italy, July 28- August*
831 *2, 2019, Volume 1: Long Papers*, pages 3721–3730.
832 Association for Computational Linguistics.

833 Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur
834 Szlam, Douwe Kiela, and Jason Weston. 2018. Per-
835 sonalizing dialogue agents: I have a dog, do you have
836 pets too? In *Proceedings of the 56th Annual Meet-*
837 *ing of the Association for Computational Linguistics,*
838 *ACL 2018, Melbourne, Australia, July 15-20, 2018,*
839 *Volume 1: Long Papers*, pages 2204–2213.

840 Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen,
841 Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing
842 Liu, and Bill Dolan. 2019b. [Dialogpt: Large-scale](#)
843 [generative pre-training for conversational response](#)
844 [generation](#). *CoRR*, abs/1911.00536.

845 Tiancheng Zhao, Ran Zhao, and Maxine Eskénazi. 2017.
846 Learning discourse-level diversity for neural dialog
847 models using conditional variational autoencoders.
848 In *Proceedings of the 55th Annual Meeting of the As-*
849 *sociation for Computational Linguistics, ACL 2017,*
850 *Vancouver, Canada, July 30 - August 4, Volume 1:*
851 *Long Papers*, pages 654–664.

Table 4: Statistics of the reddit and weibo datasets.

	Reddit	Weibo
# Users	78,031	46,973
Avg. history length	72.385	30.830
Avg. # words of query	19.8	22.9
Avg. # words of response	9.1	9.6
# Training samples	5,734,129	1,495,149
# Validation samples	10,000	10,000
# Testing samples	10,000	10,000

Table 5: The results of extra evaluation metrics on Weibo dataset and Reddit Dataset. “†” indicates that our model achieves significant improvement in t-test with p -value < 0.05 . The best results are in bold.

	BLEU-3	BLEU-4	ROUGE-1	ROUGE-2	
Weibo	Seq2Seq	0.011†	0.001†	8.740†	0.373†
	MMI	0.046†	0.007†	5.316†	0.105†
	DialoGPT	0.114†	0.027†	9.414†	0.632†
	GPMN	0.359†	0.066†	8.086†	0.753†
	PerCVAE	0.466†	0.089†	7.946†	0.485†
	Speaker	0.107†	0.041†	7.997†	0.155†
	PersonaWAE	0.889†	0.155†	11.341†	0.358†
	DHAP	1.170†	0.401†	14.131†	3.608†
	MSP (Ours)	3.973	3.522	16.249	5.812
Reddit	Seq2Seq	0.007†	0.001†	3.989†	0.233†
	MMI	0.007†	0.003†	3.960†	0.245†
	DialoGPT	0.054†	0.010†	8.977†	0.610†
	GPMN	0.039†	0.006†	4.896†	0.330†
	PerCVAE	0.068†	0.009†	8.004†	0.540†
	Speaker	0.021†	0.005†	4.017†	0.245†
	PersonaWAE	0.029†	0.007†	8.247†	0.517†
	DHAP	0.079†	0.013†	10.680†	0.697†
	MSP (Ours)	0.106	0.019	11.078	0.745

A Statistics of The Datasets

We adopt our experiment on two different datasets. The Weibo dataset contains about 46K users and about 1.5M samples (query, response, dialogue history), and the Reddit dataset has 78K users and 5.7M samples. The details of statistical information are shown in Table 4.

B Implement Details

We experiment with multiple sets of parameters to select the best model, and the final parameters are as follows: The dimension of the embeddings and Transformer hidden units is 768. The number of heads in the Transformer is 12. We use 12 layers in the decoder and 2 layers in the query encoder. The topic number is 15, and the similar user number is set to 10. The selected profile token number is 200 for the Weibo dataset and 30 for the Reddit dataset. The batch size is 128. Following (Holtzman et al., 2020), we adopt nucleus sampling as our

Table 6: Criteria of human evaluation.

Metric	Score	Criteria
Read.	1	not a complete sentence or hard to read
	2	grammatically formed
	3	fluent and well to read
Infor.	1	meaningfulness sentence
	2	contains few informative words
	3	have a clear and specific meaning
Per.	0	doesn’t resemble any user history
	1	reflect some personal information as same as user history

decoding strategy. We use the Adam (Kingma and Ba, 2015) optimizer for refiner and AdamW with a warm-up method for the generator to optimize the parameters in our model and adopt the suggested hyper-parameters for optimization.

C Extra Experimental Results

As n-gram word overlap metrics can reflect user speaking style more accurately, we evaluate the BLEU-3/4 (Papineni et al., 2002), and the result is shown in Table 5. It is consistent with other evaluations that our model outperforms every indicator. This demonstrates that user profiles also contain speaking style information, and our model can use the information to achieve a personalized response.

D Criteria of Human Evaluation

Following (Chan et al., 2019), we adopt three aspects to evaluate the generated response. *i.e.*, readability: is the response grammatically formed and smooth; informativeness: does the response contains informative words; personalization: does the response resembles any user history. The details of scoring criteria are shown in Table 6.

E Case Study

To show the effect of our model more concretely, we adopt a case study, and the results are shown in Table 7. It shows that our model can extract profiles from both current and similar users and generate informative and personalized responses. Specifically, in his dialogue history, he mentioned sports **H1**, **H2** and music **H3**, **H4** topics. Firstly, we can select similar users who also talk about sports and music, using the user refiner. Then, as the query is related to music topic, we extract the current and sim users’ dialogue history responses

Table 7: Case study for one example. Due to space limitations, we present a few user history responses and few references.

History	H1: Liverpool’s configuration has the life of a champion, support it! H2: Champion Liverpool! H3: I like to listen to music. H4: The songs on this album are really beautiful, so worth enjoying.
Query	New album "How Am I? -The sun rises as it always does", the first single "Ember" heals the system, go listen!
Persona Reference	R1: I like to listen to music. R2: The songs on this album are really beautiful, so worth enjoying.
Sim Reference	R1: Quite like this type of song. R2: Angela Leung’s "Ember" is really good.
Persona Profile	like, song, album, beautiful, enjoying
Sim Profile	like, song, Angela Leung, good, Ember
Response	DialoGPT: My little heart has flown. MSP(Ours): Angela Leung’s song is very good. Golden: Angela Leung’s songs must be listened to.

905 about music topic **R1**, **R2** by topic refiner. Further-
906 more, the token refiner selects some meaningful
907 and personalized words from long sentences of
908 reference. In this case, we can find that the token
909 refiner extracts some compliment words (like, beau-
910 tiful, enjoying) from current user’s history sentence
911 since the current user likes listening to music. And
912 the token refiner captures more concrete tokens
913 from sim users’ history sentence, such as “Angela
914 Leung” and “Ember”. By combing two profiles,
915 our personalized generator gets an informative and
916 personalized response close to ground-truth. In
917 contrast, DialoGPT generates a fluent but meaning-
918 less response to the query.